# Architectural constraints required to attain 1 Exaflop/s for scientific applications

Abhinav Bhatele, Pritish Jetley, Hormozd Gahvari,
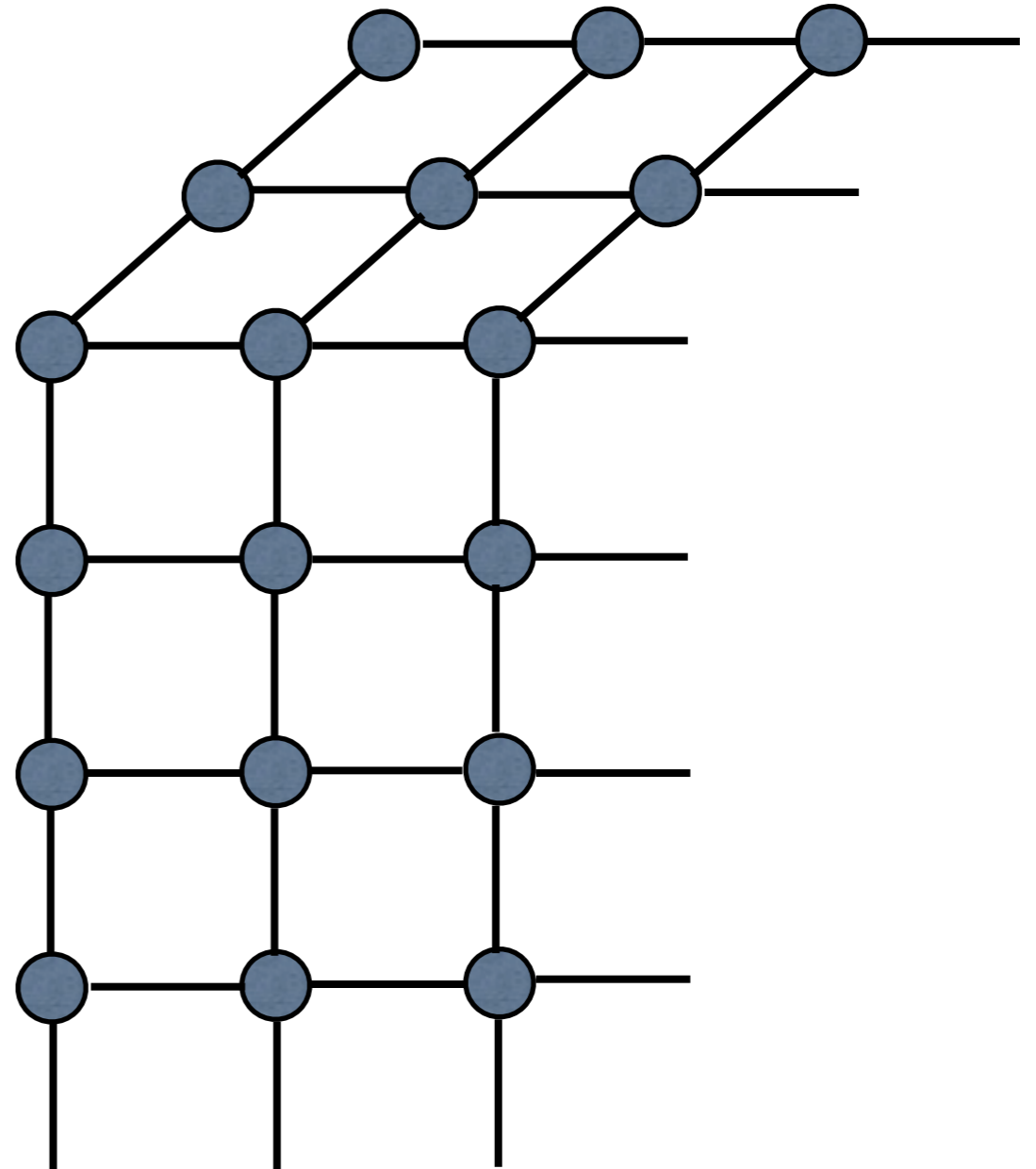Lukasz Wesolowski, William D. Gropp, Laxmikant V. Kale

Department of Computer Science
University of Illinois at Urbana-Champaign

# Motivation

- First Teraflop/s computer (ASCI Red, 1997), first Petaflop/s computer (RoadRunner, 2008), Exaflop/s 2018 ?

- Hardware challenges: power/energy, memory, communication

- Software challenges: algorithms and implementations that will scale

- Architectural features to attain 1 Exaflop/s ?

# A possible exascale machine

- $2^{20} = 1,048,576$ nodes

- $2^{10}$ cores per node

- 10 Gflop/s cores, time to compute a flop, $t_c = 0.1$ ns

- 10.74 Exaflop/s peak performance

# Modeling methodology

- Estimate the floating point calculations/operations per iteration,

$$T_{comp} = \frac{1}{\eta} \times f(N, P_c) \times n \times t_c$$

- Time for communication based on number and size of messages

$$T_{comm} = M \times (t_s + h(N, P_c) \times t_w)$$

- Using total number of floating point operations and time per iteration, $\frac{flops}{T} > 10^{18}$

# Applications

- Molecular Dynamics

  - Short-range forces, spatial decomposition

- Cosmological Simulations

  - Tree algorithms

- Unstructured grid problems

  - Finite element solvers

PPL
UIUC

# Molecular Dynamics

- ## Spatial decomposition

---

**Algorithm 1** Computation in one time step of MD

---

Receive atoms from neighboring processors
**for** $i = 1$ to $N_p$ **do**
    **for** $j = 1$ to $N_i$ **do**
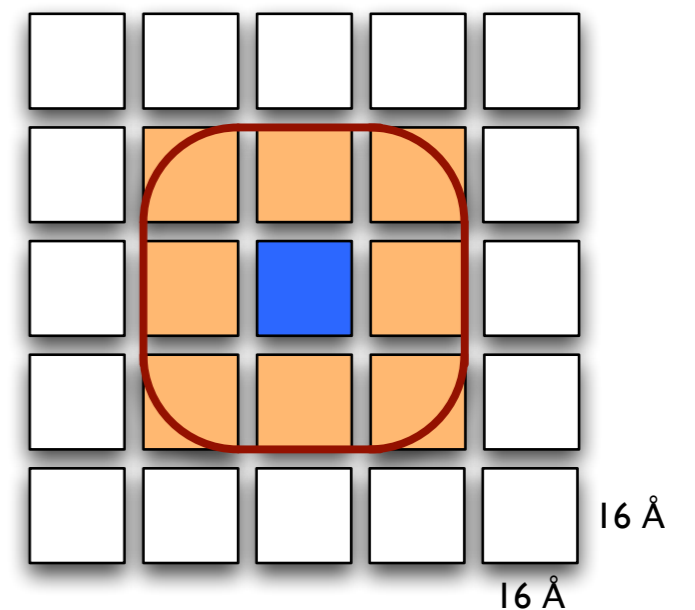        **if** atoms are within cutoff radius, $r_c$ **then**
            Compute forces on pairs of atoms
        **end if**
    **end for**
**end for**
Update atom positions and velocities

---

16 Å

16 Å

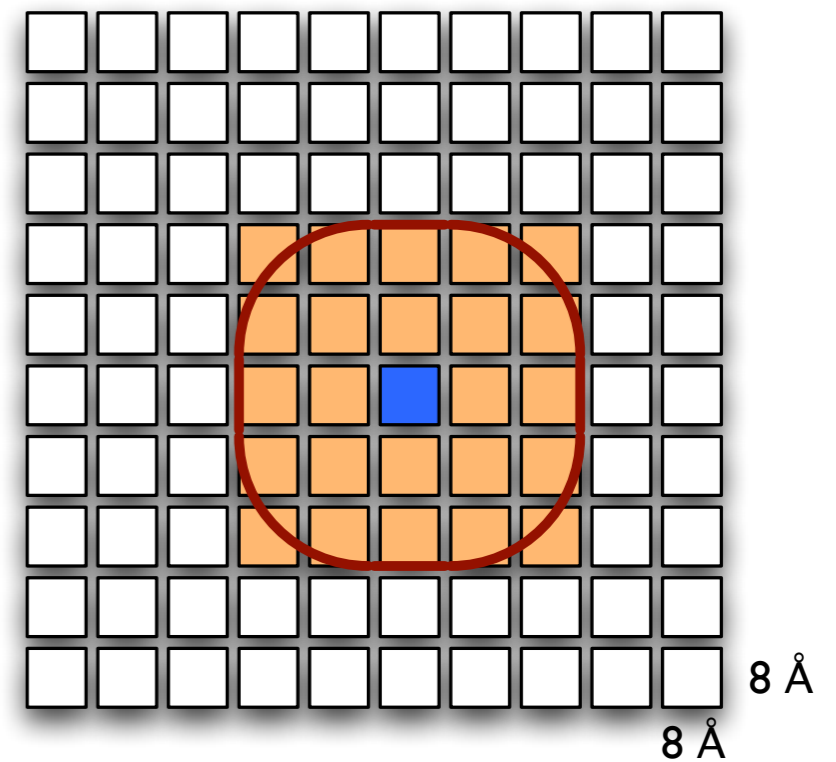PPL
UIUC

# Weak scaling of MD

- Size of molecular system = $100 * 2^{30}$ = 107 billion atoms

- Number of floating point operations = $33547 * N$

$$\frac{flops}{T} > 10^{18}$$

$$\frac{33547 \times N}{10^{18}} > T$$

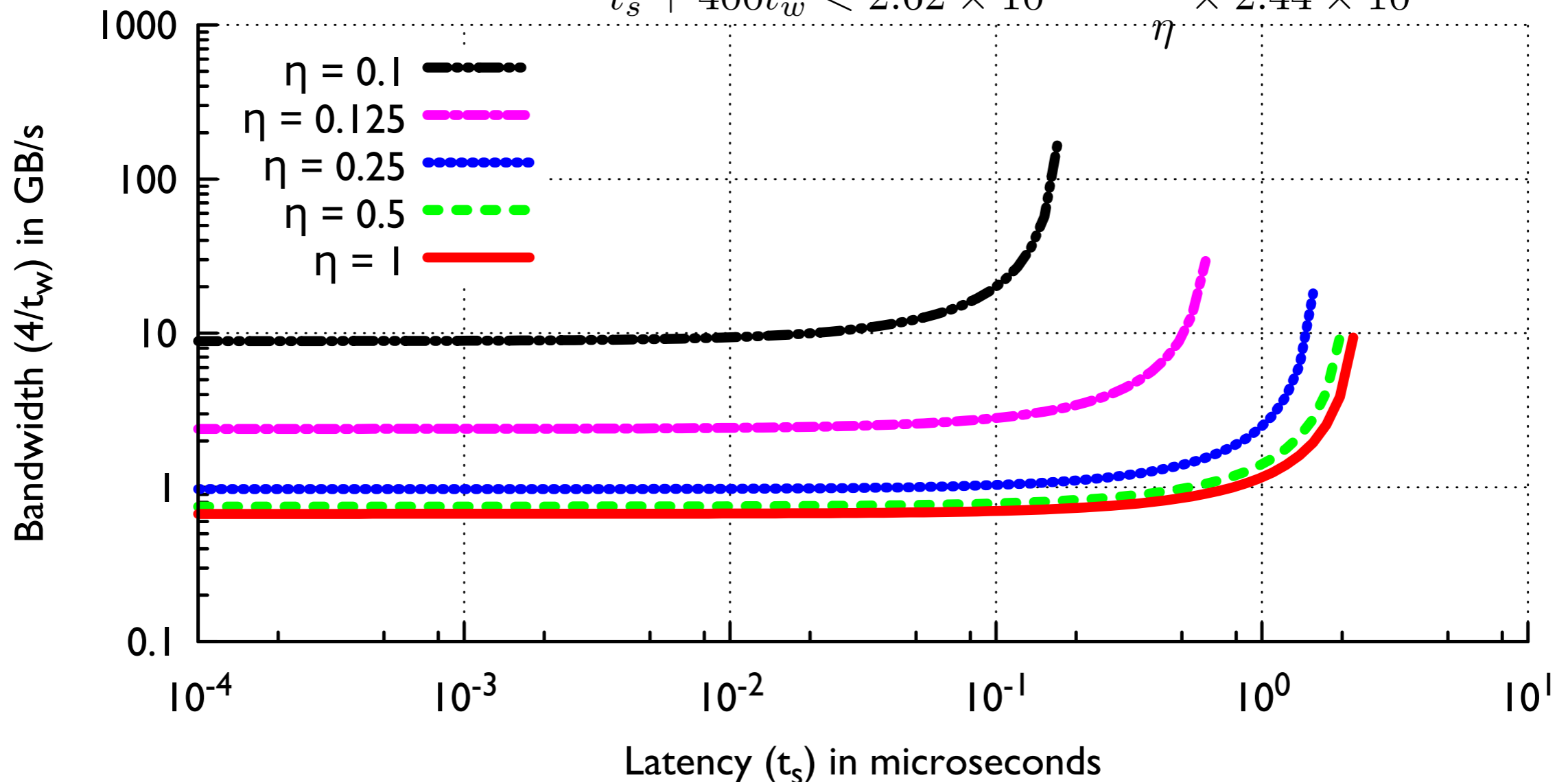- Putting N = $100 * 2^{30}$,

$$T < 3.6 \times 10^{-3}$$

PPL
UIUC

- 100 atoms per cell

- Split the cells in two of the three dimensions

- Each cell communicates with 5*5*3 = 75 other cells

- For a block of 8*8*16 cells placed on a node only the ones on the boundary communicate inter-node



8 Å

8 Å

PPL

UIUC

1000

8

η = 0.125

# Inferring network parameters

$$\frac{1}{\eta} \times \frac{N}{P_c} \times 33547 \times t_c + 1376 \times \left(t_s + \frac{N}{P_c} 4t_w\right) < 3.6 \times 10^{-3}$$

$$t_s + 400t_w < 2.62 \times 10^{-6} - \frac{1}{\eta} \times 2.44 \times 10^{-7}$$

# Inferring network parameters

$$\frac{1}{\eta} \times \frac{N}{P_c} \times 33547 \times t_c + 1376 \times \left( t_s + \frac{N}{P_c} 4t_w \right) < 3.6 \times 10^{-3}$$
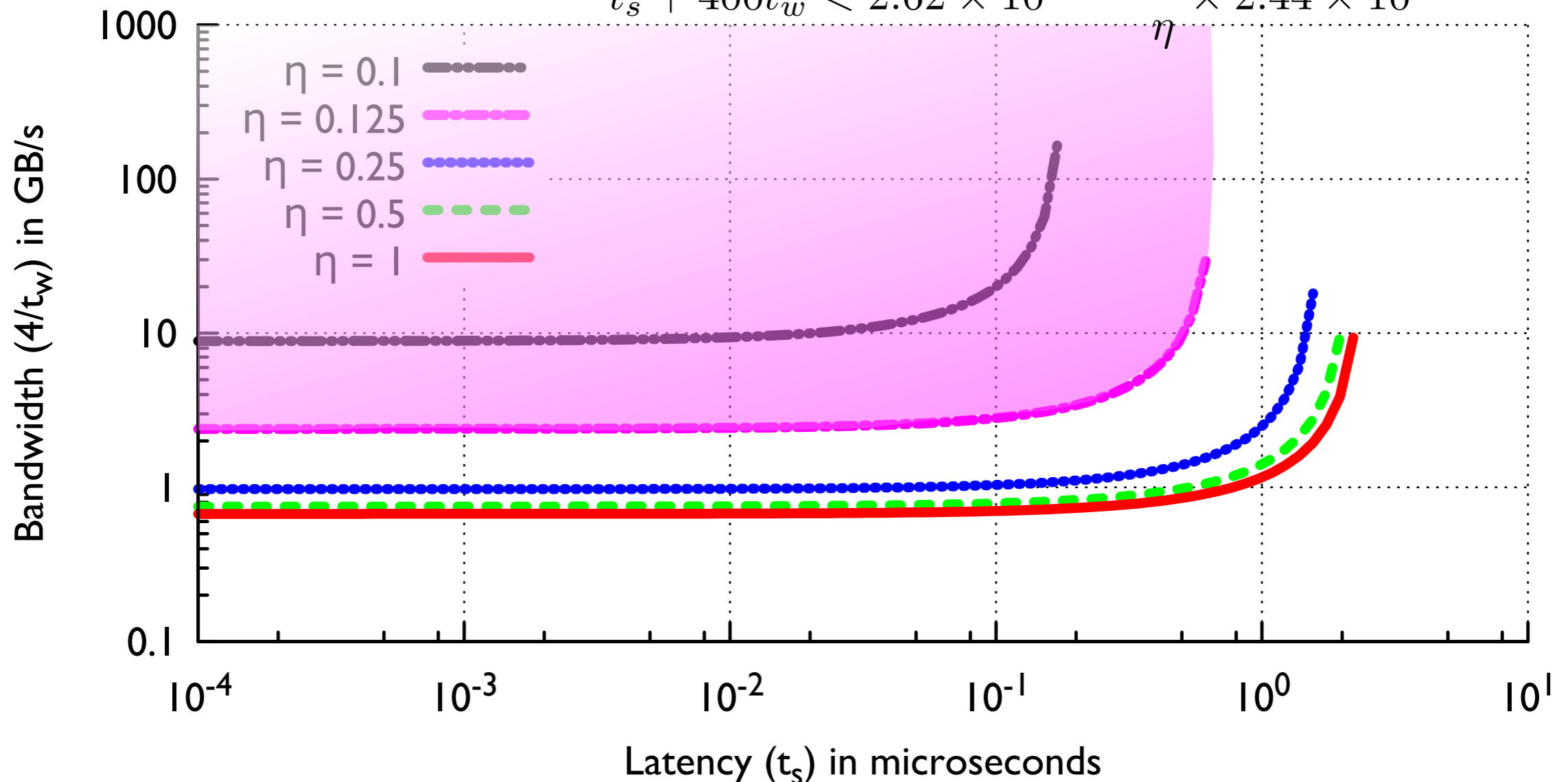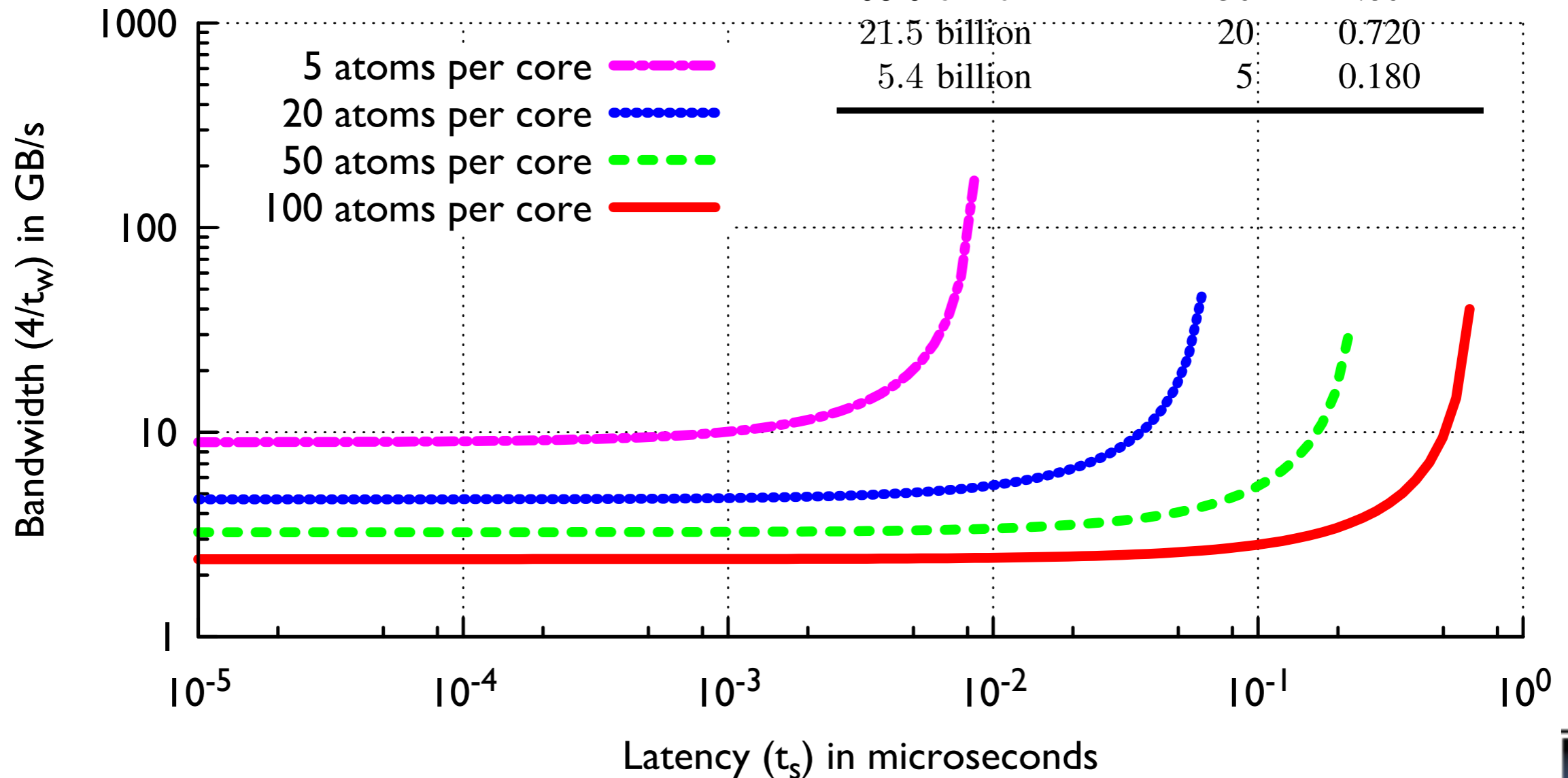
$$t_s + 400t_w < 2.62 \times 10^{-6} - \frac{1}{\eta} \times 2.44 \times 10^{-7}$$



Bandwidth $(4/t_w)$ in GB/s

- η = 0.1
- η = 0.125
- η = 0.25
- η = 0.5
- η = 1

Latency $(t_s)$ in microseconds

# Smaller problem sizes

| # Atoms | Atoms/core | Time (ms) |
|---|---|---|
| 107 billion | 100 | 3.602 |
| 53.6 billion | 50 | 1.801 |
| 21.5 billion | 20 | 0.720 |
| 5.4 billion | 5 | 0.180 |



Legend:
- 5 atoms per core
- 20 atoms per core
- 50 atoms per core
- 100 atoms per core

X-axis: Latency ($t_s$) in microseconds
Y-axis: Bandwidth ($4/t_w$) in GB/s

PPL
UIUC

# Computational Cosmology

- Several approaches to computing trajectories of bodies under gravitational attraction

  - Direct, all-pairs

  - Tree-based approximate methods

  - Particle-mesh or "grid" methods

- We consider locality-aware tree codes

# Modeling problem size

- What problems will be of interest given an exascale-level machine

- Extrapolate from current state-of-the-art simulations

- About 8192 particles are required per core for good parallel efficiency at petascale

- Given O(N log N)/P work per core, about 6350 particles per core are needed at exascale (total 6.8 trillion)

PPL
UIUC

# Barnes-Hut computation

- Analyze algorithm:

  - Domain decomposition => distributed spatial tree

  - Every processor core gets a number of leaves

  - For each leaf l, Traverse(l, root)

```
Traverse(leaf l, node n) {
  if(IsLeaf(n)) {
    LeafForces(l, n);
  }
  else if(Side(n)/|r(n)-r(l)| < Θt)
{
    CellForces(l, n);
  }
  else {
    foreach(node c in Children(n)) {
      Traverse(l, c);
    }
  }
}
```

# Total computation

- Number of floating point operations per iteration,

$$312 \times 77 \times N \times \lg \frac{N}{B} + 38 \times 33 \times B \times N$$

- To attain a rate of 1 Exaflop/s,

$$\frac{24024 \times N \lg(N/B) + 1254 \times BN}{T} > 10^{18}$$

- T < 6.52s

# Total communication

IPDPS 2011 © Abhinav Bhatele

# Total communication

- Could obtain communication from number of expansions E(l) for every level l

# Total communication

- Could obtain communication from number of expansions $E(l)$ for every level $l$

- However, cores on an SMP node can reuse remote data through software caching

PPL
UIUC

# Total communication

- Could obtain communication from number of expansions E(l) for every level l

- However, cores on an SMP node can reuse remote data through software caching

- Communication with remote data caching:

  - Each SMP node holds a cube of space

  - Cores holding particles near surface of cube request remote data - other cores reuse data

  - Find each SMP node's *halo* of requests at each level of tree

# Communication analysis

Leaf level:
$$12n_b^2 + 36n_b + 8$$

1 level above leaves:
$$12(n_b/2)^2 + 36(n_b/2) + 8$$

2 levels above leaves:
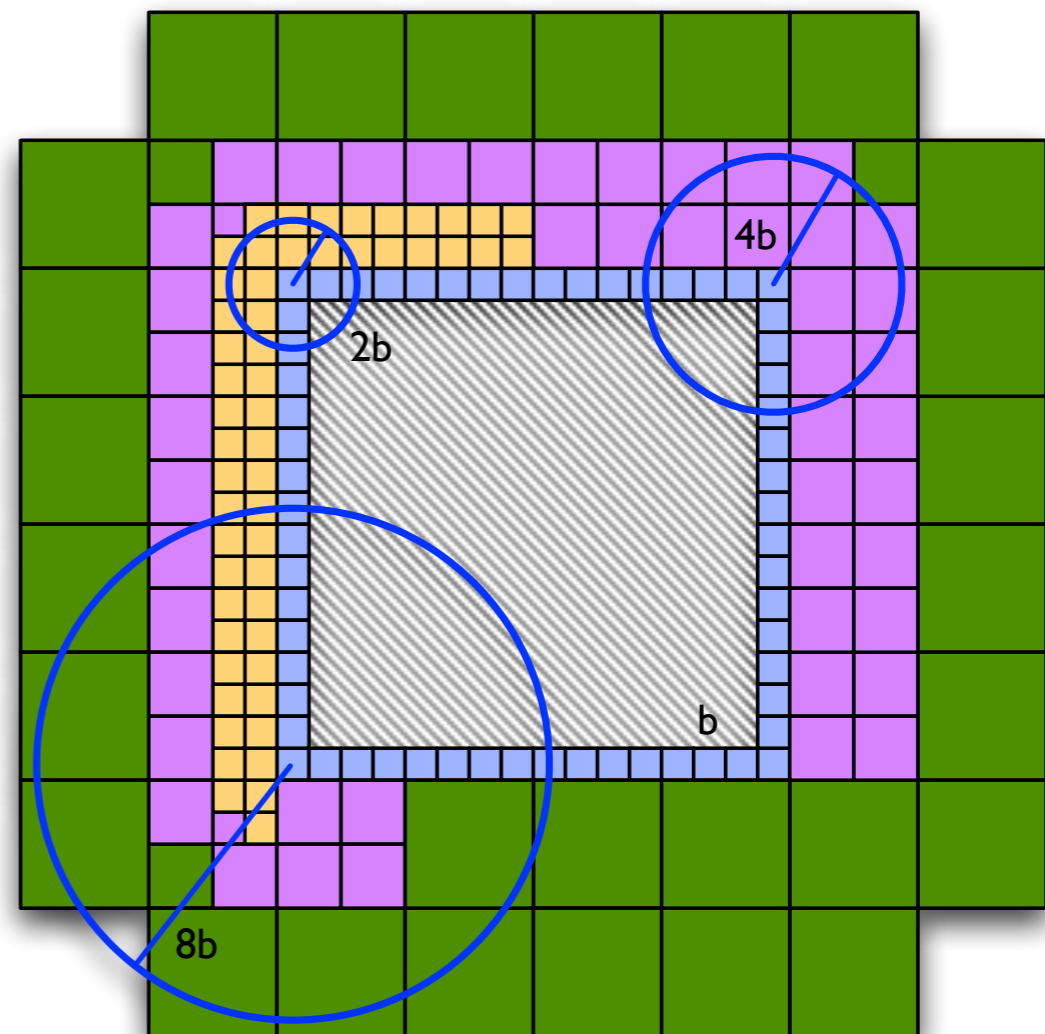$$12(n_b/4)^2 + 36(n_b/4) + 8$$

3 levels above leaves:
$$12(n_b/8)^2 + 36(n_b/8) + 8$$

…

Total:

$$C_1^{\text{cell}} = \sum_{i=0}^{\lg n_b} \left( 12 \left( \frac{n_b}{2^i} \right)^2 + 36 \left( \frac{n_b}{2^i} \right) + 8 \right)$$
$$= 16n_b^2 + 72n_b + 8 \lg n_b - 32 \quad \text{cells}$$

# Upper-level calls

- Previous reasoning valid as long as edge length of requested calls $<= c/(P_n)^{1/3}$

- Use reasoning similar to calculation of E(l) to get number of larger, upper-level cells requested per SMP node,
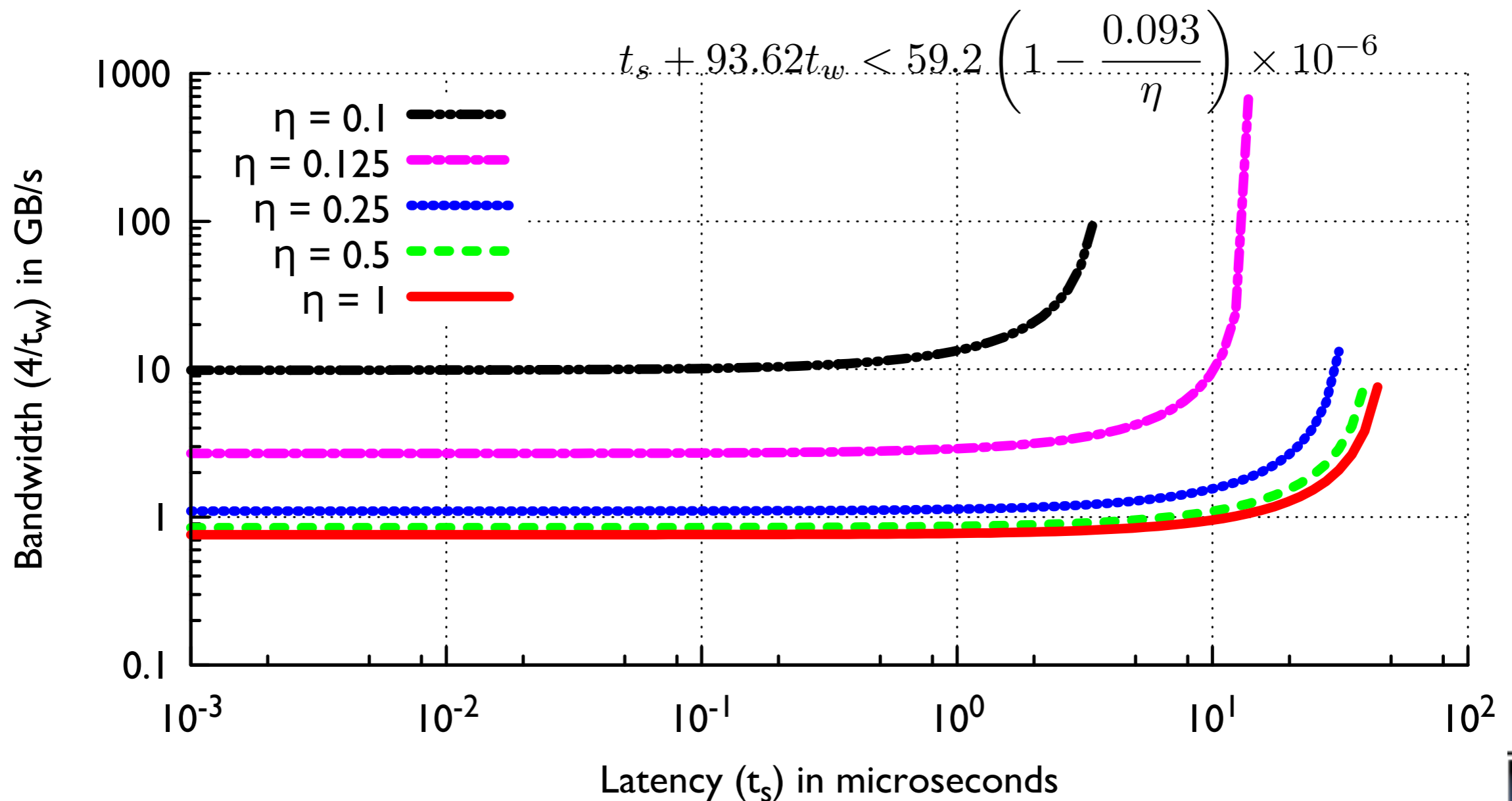
- Previous reasoning valid as long as edge length of requested calls <= c/(P$_n$)$^{1/3}$

- Use reasoning similar to calculation of E(l) to get number of larger, upper-level cells requested per SMP node,

$$C_2^{\text{cell}} = 31 \left( \frac{\lg P_n}{3} - 1 \right) \quad \text{cells}$$
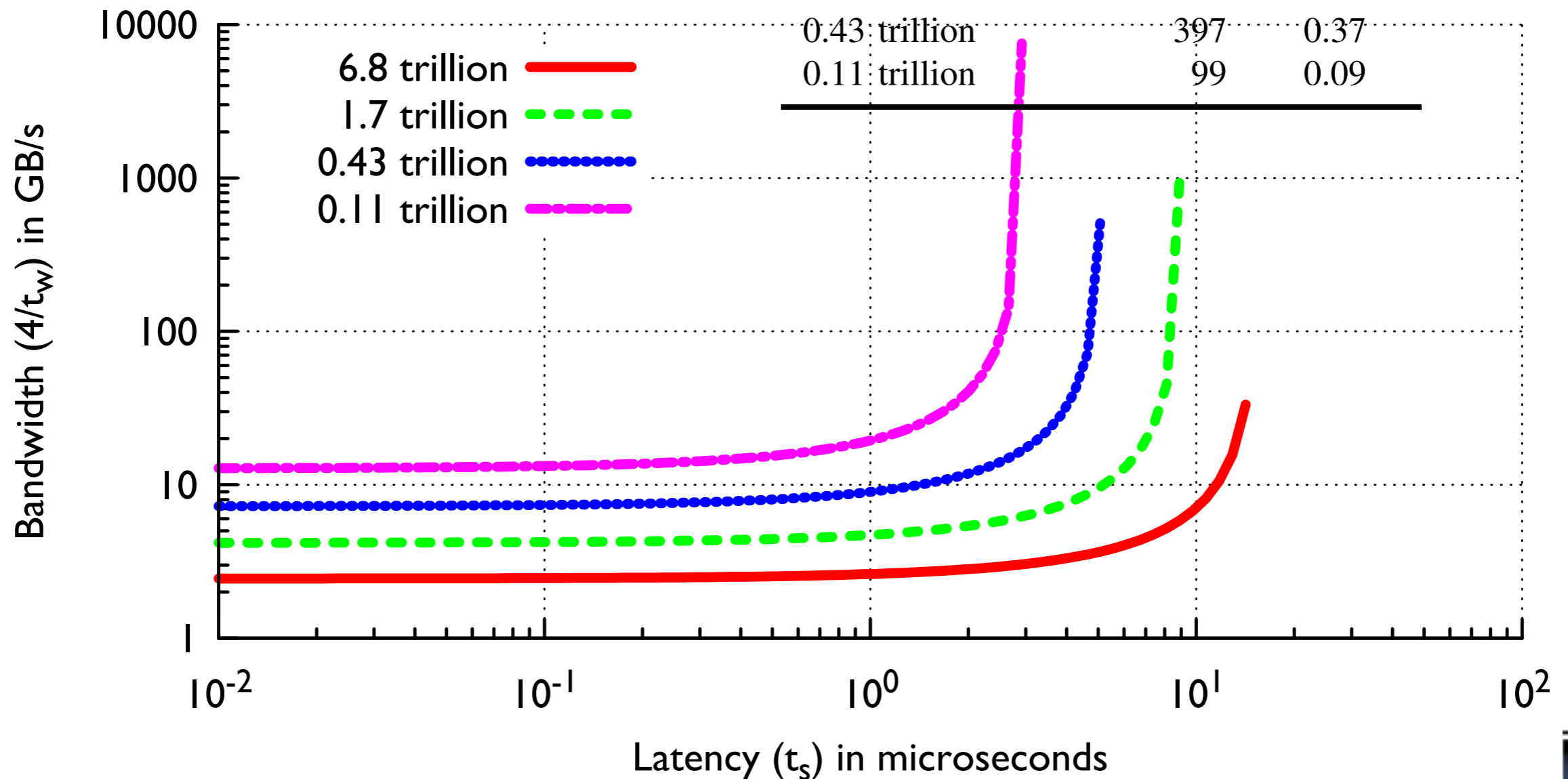
$$T_{comm} = 15946(t_s + 56t_w) + 93968(t_s + 100t_w)$$

# Inferring network parameters

$$\frac{6.52 \times 10^{18}}{P_c} \times \frac{t_c}{\eta} + (1.1 \times 10^5 t_s + 1.03 \times 10^7 t_w) < 6.52$$

$$t_s + 93.62 t_w < 59.2 \left(1 - \frac{0.093}{\eta}\right) \times 10^{-6}$$

# Smaller problem sizes

| # Particles | Particles/core | Time (s) |
|---|---|---|
| 6.8 trillion | 6350 | 6.52 |
| 1.7 trillion | 1588 | 1.55 |
| 0.43 trillion | 397 | 0.37 |
| 0.11 trillion | 99 | 0.09 |



Legend:
- 6.8 trillion
- 1.7 trillion
- 0.43 trillion
- 0.11 trillion

Y-axis: Bandwidth $(4/t_w)$ in GB/s

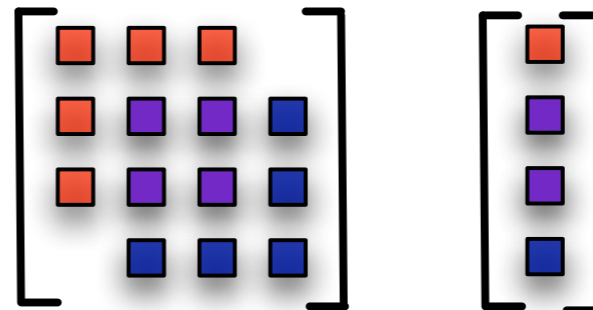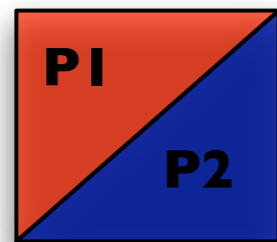X-axis: Latency $(t_s)$ in microseconds

# Finite Element Solvers

- Method of choice for unstructured grid problems

- Involves two phases:

  - Assembly: put linear system together

  - Solve: the system

- Linear problems: one assembly, one (time-independent) or more (time-dependent) solves

- Nonlinear problems: repeat assembly/solve process until convergence

# Approach to Solution

- Based on recent work by Sahni et al. that scales FEM to near-petascale

- Partition the problem by elements, storing shared DOFs redundantly



- Assembly becomes nearest-neighbor: focus on solve

# Approach to Solution

- ## Assume conjugate gradient linear solver

  - Setup: one mat-vec product, one vector subtraction, one dot product

  - Iteration loop: one mat-vec product, two vector additions, one vector subtraction, two dot products

---

**Algorithm 3** $\text{CG}(A, b, x_0, rtol)$

---

$r_0 \leftarrow b - Ax_0$

$p_0 \leftarrow r_0$

$k \leftarrow 0$

**while** $||r_k||_2 \geq rtol$ **do**

$\quad \alpha_k \leftarrow \frac{r_k^T r_k}{p_k^T A p_k}$

$\quad x_{k+1} \leftarrow x_k + \alpha_k p_k$

$\quad r_{k+1} \leftarrow r_k - \alpha_k A p_k$

$\quad \beta_k \leftarrow \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$

$\quad p_{k+1} \leftarrow r_{k+1} + \beta_k p_k$
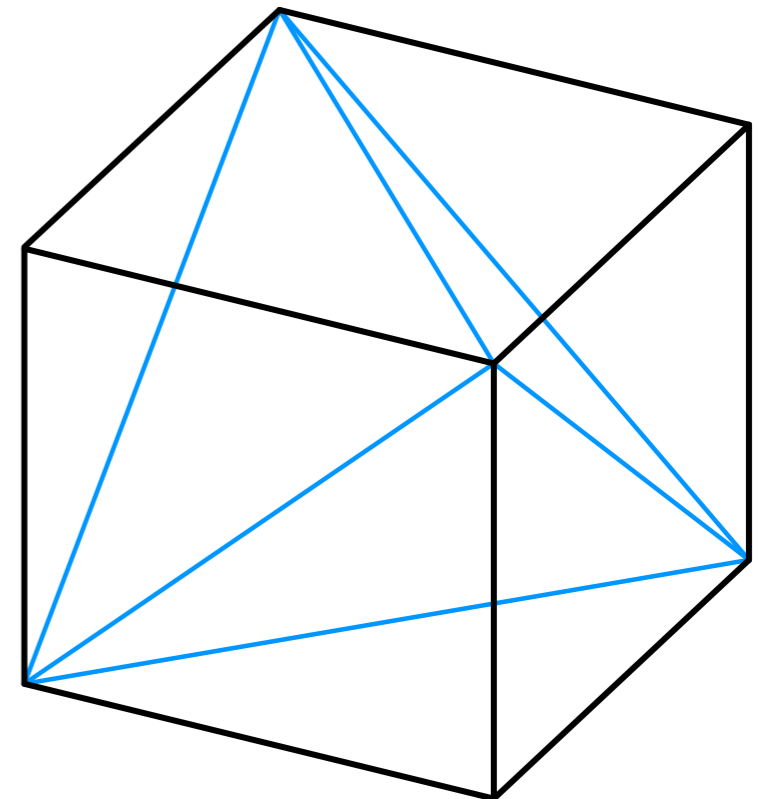
$\quad k \leftarrow k + 1$

**end while**

return $x_k$
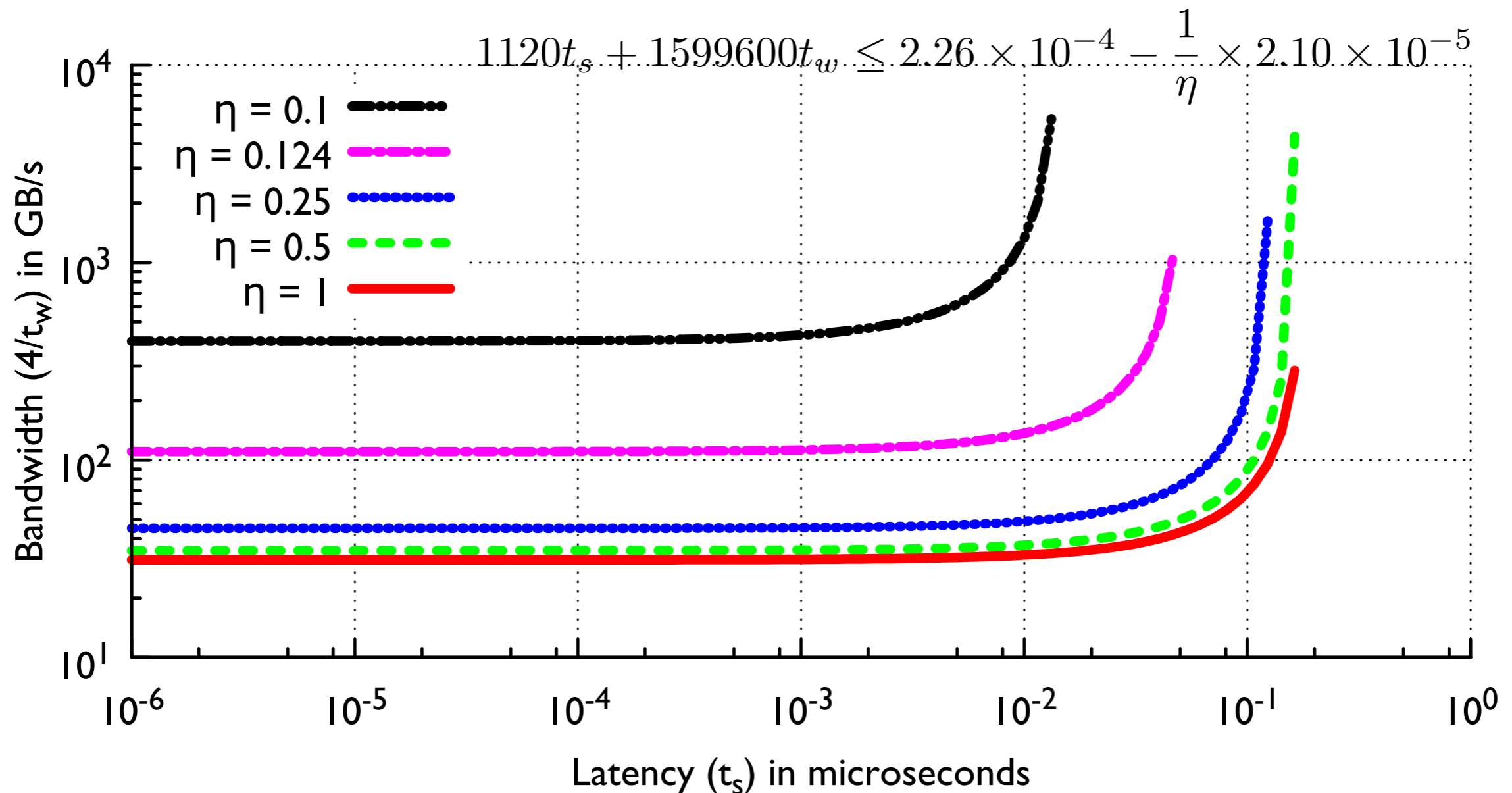
---

# FEM: Weak Scaling

- Consider problem on 3D cubic tet mesh

  - Each core gets $16^3$ cubes

  - Degrees of freedom on each processor = $17^3$

- Global DOFs = 4.4 trillion

- Solve time per iteration:

$$T_{CG}^{iter} = \frac{1}{\eta} \times \left( (2s_i + 6)n_i + \frac{N}{P_c} + 2\lg P_n \right) t_c$$
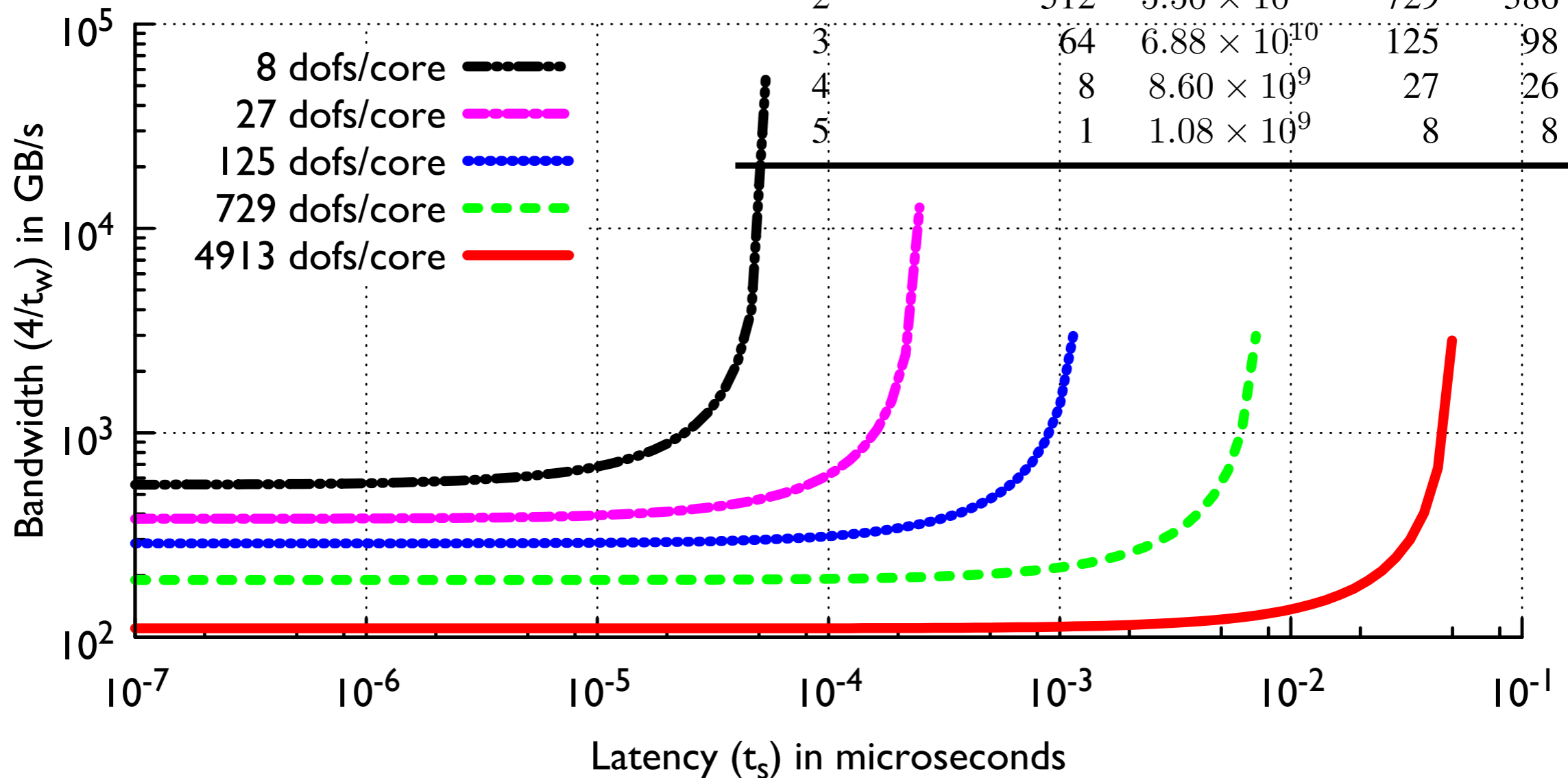$$+ 2(520 + 2\lg P_n)t_s$$
$$+ 2(520\tilde{n}_i + 2\lg P_n)t_w$$

# Inferring network parameters

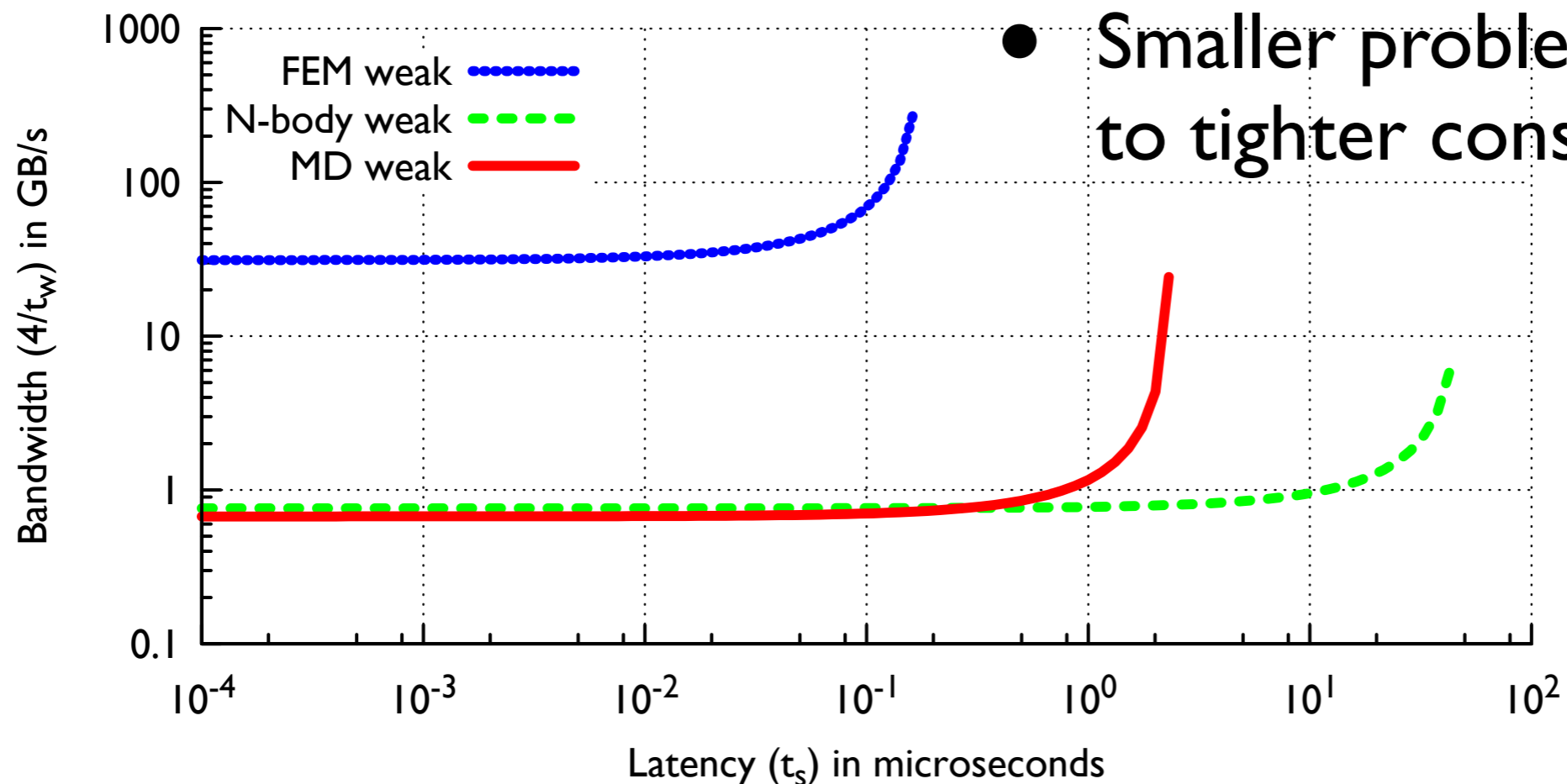$$\frac{P_c(2s_i + 1)n_i + N + 2 \lg P_n}{T_{\mathrm{CG}}^{\mathrm{iter}}} \geq 10^{18}$$

$$1120 t_s + 1599600 t_w \leq 2.26 \times 10^{-4} - \frac{1}{\eta} \times 2.10 \times 10^{-5}$$



Feasibility Region for CG

Bandwidth ($4/t_w$) in GB/s

Latency ($t_s$) in microseconds

# Smaller problem sizes

| Problem | Cubes/core | $N$ | $n_i$ | $\tilde{n}_i$ |
|---|---|---|---|---|
| 1 | 4096 | $4.40 \times 10^{12}$ | 4913 | 1538 |
| 2 | 512 | $5.50 \times 10^{11}$ | 729 | 386 |
| 3 | 64 | $6.88 \times 10^{10}$ | 125 | 98 |
| 4 | 8 | $8.60 \times 10^{9}$ | 27 | 26 |
| 5 | 1 | $1.08 \times 10^{9}$ | 8 | 8 |



Legend: 8 dofs/core, 27 dofs/core, 125 dofs/core, 729 dofs/core, 4913 dofs/core

Y-axis: Bandwidth ($4/t_w$) in GB/s

X-axis: Latency ($t_s$) in microseconds

# Summary

- Modest communication requirements for MD and cosmology at exascale

- Smaller problem sizes lead to tighter constraints

# Future work

- Research required in area of communication-minimizing algorithms and high-bandwidth low-latency networks

- Detailed analysis of each application class

  - MD: long-range forces

  - Cosmology: particle-mesh methods

  - FEM: other solvers, preconditioning

- Studies for specific networks and contention